

Here I define the Bayesian model I want to use to test the performance of the different priors. The features I want to include are the following: continuous outcome, any number/measurement scale.

Let's start with the **model**

$$\begin{aligned} y_{ij} &= \mathbf{x}_{ij}^T \boldsymbol{\theta} + \mathbf{z}_{ij}^T \mathbf{b}_i + \epsilon_{ij} \\ \mathbf{b}_i &\sim N(\mathbf{0}, \boldsymbol{\Psi}) \\ \epsilon_{ij} &\sim N(0, \sigma^2) \end{aligned} \tag{1}$$

(Vectors are in bold, matrix are capital Greek letters).

I'm going to define the following **priors**:

$$p(\boldsymbol{\theta}) \propto 1 \tag{2}$$

$$p(\sigma^2) \propto \sigma^{-2} \tag{3}$$

For what concerns the random effects variance covariance matrix, different priors are tested. In particular we used:

- inverse-Wishart

$$p(\boldsymbol{\Psi}) \propto IW(\nu, S_0) \tag{4}$$

where we choose $\nu = k - 1 + e$, and $S_0 = \text{diag}(k - 1 + e)$, following indications by Gelman *et al* (2014). Thanks to the partitioning property, we know that the diagonal components of $\boldsymbol{\Psi}$ have an inverse-Wishart distribution themselves. Furthermore, for a univariate case, we know that an inverse Wishart distribution simplifies to an inverse Gamma with parameters $k = 1, \alpha = \frac{\nu}{2}, \beta = \frac{S_{0kk}}{2}$. The inverse-Wishart priors we are trying are:

Prior Description	ν	S_0
1. IW uninformative	$k - 1 + e$	$\text{diag}(k-1+e)$
2. IW educated	$k - 1 + e$	educated guess

- inverse-Wishart *a là* Huang and Wand

$$\begin{aligned} p(\boldsymbol{\Psi}|a_1, a_2) &\propto IW(\nu + k - 1, 2\nu \times \text{diag}(1/a_1, 1/a_2)), \\ a_k &\propto IG(1/2, 1/A_k^2), \end{aligned} \tag{5}$$

with $\nu = 2$ and $\mathbf{A} = [1000, 1000]$. The marginal distribution of any standard deviation term in $\boldsymbol{\Psi}$ is Half- $t(\nu, A_k)$ and, when choosing $\nu = 2$, the marginal distribution on the correlation term is uniform on $(-1, 1)$, see property 2 to 4 in Huang and Wand (2013, p. 442). Furthermore, according to Huang and Wand (2013, p. 441) arbitrarily large values for a_k lead to arbitrarily weak priors on the standard deviation term. Hence, our choices for the parameter of this prior are:

Prior Description	ν	\mathbf{A}
3. IW <i>a là</i> HW	2	$[1000, 1000]$

- Matrix-F variate

$$\begin{aligned} p(\boldsymbol{\Psi}) &\propto F(\boldsymbol{\Psi}; \nu, \delta, \mathbf{B}) \\ &\propto \int IW(\boldsymbol{\Psi}; \delta + k - 1, \Sigma) \times W(\boldsymbol{\Sigma}; \nu, \mathbf{B}) d\boldsymbol{\Sigma} \end{aligned} \tag{6}$$

with degrees of freedom $\nu > k - 1$, $\delta > 0$, and \mathbf{B} a positive definite scale matrix that functions as prior guess. Three different choices were made for \mathbf{B} in this paper: $\text{diag}(10^3)$, proper neighbor of $(\sigma^2)^{-\frac{1}{2}}$; \mathbf{B}_{ed} , an educated guess based on data exploration, \mathbf{R}^* and an empirical bayes choice following Kass and Natarajan (2006).

Considering a 2×2 random effects variance covariance matrix (random intercepts, and random slopes) that is matrix-F distributed, $F(\nu, \delta, \mathbf{B})$, the marginal distribution on the standard deviations of the random effects are univariate $F(\nu, \delta, b_{11})$ and $F(\nu, \delta, b_{22})$, with $\nu > 1, \delta > 0, b_{jj} > 0$. There we chose the first integer number we could for the parameters ν , and δ . When I chose $\delta = 2$, I was thinking of staying close to what has been done in Mulder Pericchi. When I chose $\delta = e$ I wanted to try with some uninformative prior that got as close as possible to non-informative.

Prior Description	ν	δ	S_0
4. mat-F proper neighbor	2	2	$10^3 \times \mathbf{I}_2$
5. mat-F uninformative	$k - 1 + e$	e	educated guess
6. mat-F educated guess	2	2	educated guess
7. mat-F empirical Bayes	2	2	\mathbf{R}^*

The derivation of the conditional posterior follows.

Full conditional for θ (fixed effects)

Let's start with

$$p(\theta | \mathbf{y}, \mathbf{X}, \mathbf{Z}, \mathbf{b}_i, \Psi, \sigma^2) = p(\mathbf{y} | \theta, \mathbf{X}, \mathbf{Z}, \mathbf{b}_i, \Psi, \sigma^2) p(\theta)$$

where

$$\begin{aligned} p(\mathbf{y} | \theta, \mathbf{X}, \mathbf{Z}, \Psi, \sigma^2) &= \prod_{i=1}^n \prod_{j=1}^J p(y_{ij} | \theta^T \mathbf{x}_{ij} + \mathbf{b}_i^T \mathbf{z}_{ij}, \Psi, \sigma^2) \\ &\propto \exp\left(-\frac{1}{2\sigma^2} SSR\right) \end{aligned}$$

and

$$SSR = \sum_{i=1}^n \sum_{j=1}^J (y_{ij} - \theta^T \mathbf{x}_{ij} - \mathbf{b}_i^T \mathbf{z}_{ij})^2$$

where can rewrite y_{ij} as \tilde{y}_{ij} , with $\tilde{y}_{ij} = y_{ij} - \mathbf{b}_i^T \mathbf{z}_{ij}$ which makes SSR:

$$\begin{aligned} SSR &= \sum_{i=1}^n \sum_{j=1}^J (\tilde{y}_{ij} - \theta^T \mathbf{x}_{ij})^2 \\ &= (\tilde{\mathbf{y}} - \mathbf{X}\theta)^T (\tilde{\mathbf{y}} - \mathbf{X}\theta) \\ &= \tilde{\mathbf{y}}^T \tilde{\mathbf{y}} - 2\theta^T \mathbf{X}\tilde{\mathbf{y}} + \theta^T \mathbf{X}^T \mathbf{X}\theta \end{aligned}$$

Hence,

$$p(\mathbf{y} | \theta, \mathbf{X}, \mathbf{Z}, \Psi, \sigma^2) \propto \exp\left(-\frac{1}{2\sigma^2} [-2\theta^T \mathbf{X}\tilde{\mathbf{y}} + \theta^T \mathbf{X}^T \mathbf{X}\theta]\right)$$

Combining this with the prior we obtain:

$$\begin{aligned} p(\theta | \mathbf{y}, \mathbf{X}, \mathbf{Z}, \Psi, \sigma^2) &\propto \exp\left(-\frac{1}{2} \theta^T \mathbf{X}^T \mathbf{X}\theta + \theta^T \mathbf{X}\tilde{\mathbf{y}}\right) \\ \theta | \cdot &\sim N\left(\frac{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}\tilde{\mathbf{y}}}{\sigma^2}, \frac{(\mathbf{X}^T \mathbf{X})^{-1}}{\sigma^2}\right) \end{aligned} \tag{7}$$

Full conditional for \mathbf{b}_i (random effects)

To derive this one we can start from:

$$p(\mathbf{b}_i|\mathbf{y}, \mathbf{X}, \mathbf{Z}, \boldsymbol{\theta}, \boldsymbol{\Psi}, \sigma^2) = p(\mathbf{y}_i|\boldsymbol{\theta}, \mathbf{b}_i, \mathbf{X}, \mathbf{Z}, \boldsymbol{\Psi}, \sigma^2)p(\mathbf{b}_i)$$

We know that

$$p(\mathbf{y}_i|\cdot) = \prod_{j=1}^J p(y_{ij}|\boldsymbol{\theta}^T \mathbf{x}_{ij} + \mathbf{b}_i^T \mathbf{z}_{ij}, \sigma^2) \propto \exp(-\frac{1}{2\sigma^2} SSR_i)$$

with

$$SSR = \sum_{j=1}^J (y_{ij} - \boldsymbol{\theta}^T \mathbf{x}_{ij} - \mathbf{b}_i^T \mathbf{z}_{ij})^2$$

and we can rewrite y_{ij} as $\tilde{y}_{ij} = y_{ij} - \boldsymbol{\theta}^T \mathbf{x}_{ij}$, which would make SSR be

$$\begin{aligned} SSR &= \sum_{j=1}^J (\tilde{y}_j - \boldsymbol{\theta}^T \mathbf{x}_j)^2 \\ &= (\tilde{\mathbf{y}} - \mathbf{b}_i^T \mathbf{Z}_i)^T (\tilde{\mathbf{y}} - \mathbf{b}_i^T \mathbf{Z}_i) \\ &= \tilde{\mathbf{y}}^T \tilde{\mathbf{y}} - 2\mathbf{b}_i^T \mathbf{Z}_j \tilde{\mathbf{y}} + \mathbf{b}_i^T \mathbf{Z}_i^T \mathbf{Z}_j \mathbf{b}_i \end{aligned}$$

Hence,

$$p(\mathbf{y}_i|\cdot) \propto \exp(-\frac{1}{2\sigma^2} [-2\mathbf{b}_i^T \mathbf{Z}_j \tilde{\mathbf{y}} + \mathbf{b}_i^T \mathbf{Z}_i^T \mathbf{Z}_j \mathbf{b}_i])$$

We also know that in this case, the "prior" is

$$p(\mathbf{b}_i) \propto N(\mathbf{0}, \boldsymbol{\Psi}) \propto \exp(-\frac{1}{2} [-2\mathbf{b}_i^T \boldsymbol{\Psi}^{-1} \mathbf{0} + \mathbf{b}_i^T \boldsymbol{\Psi}^{-1} \mathbf{b}_i])$$

In conclusion, combining the sampling model and the prior, we get:

$$\begin{aligned} p(\mathbf{b}_i|\cdot) &\propto \exp(-\frac{1}{2\sigma^2} [-2\mathbf{b}_i^T \mathbf{Z}_j \tilde{\mathbf{y}} + \mathbf{b}_i^T \mathbf{Z}_i^T \mathbf{Z}_j \mathbf{b}_i] - \frac{1}{2\sigma^2} [-2\mathbf{b}_i^T \mathbf{Z}_j \tilde{\mathbf{y}} + \mathbf{b}_i^T \mathbf{Z}_i^T \mathbf{Z}_j \mathbf{b}_i]) \\ \mathbf{b}_i| \cdot &\propto N\left(\left(\boldsymbol{\Psi}^{-1} + \frac{\mathbf{Z}_i^T \mathbf{Z}_i}{\sigma^2}\right)^{-1} \left(\boldsymbol{\Psi}^{-1} \mathbf{0} + \frac{\mathbf{Z}_i^T \tilde{\mathbf{y}}_i}{\sigma^2}\right), \left(\boldsymbol{\Psi}^{-1} + \frac{\mathbf{Z}_i^T \mathbf{Z}_i}{\sigma^2}\right)^{-1}\right) \end{aligned} \quad (8)$$

Full conditional for σ^2 (error variance)

The full conditional posterior can be expressed as:

$$p(\sigma^2|\mathbf{y}, \mathbf{X}, \mathbf{Z}, \boldsymbol{\theta}, \mathbf{b}_i, \boldsymbol{\Psi}) = p(\mathbf{y}|\boldsymbol{\theta}, \mathbf{b}_i, \mathbf{X}, \mathbf{Z}, \boldsymbol{\Psi}, \sigma^2)p(\sigma^2)$$

The sampling model is the same we saw for the full conditional distribution of $\boldsymbol{\theta}$:

$$\begin{aligned} p(\mathbf{y}|\boldsymbol{\theta}, \mathbf{X}, \mathbf{Z}, \boldsymbol{\Psi}, \sigma^2) &= \prod_{i=1}^n \prod_{j=1}^J p(y_{ij}|\boldsymbol{\theta}^T \mathbf{x}_{ij} + \mathbf{b}_i^T \mathbf{z}_{ij}, \boldsymbol{\Psi}, \sigma^2) \\ &= \prod_{i=1}^n \prod_{j=1}^J (2\pi\sigma^{-2})^{-\frac{1}{2}} \exp(-\frac{(y_{ij} - \boldsymbol{\theta}^T \mathbf{x}_{ij} - \mathbf{b}_i^T \mathbf{z}_{ij})^2}{2\sigma^2}) \end{aligned}$$

However, we are now interested in σ^2 , hence

$$\begin{aligned} p(\mathbf{y}|\boldsymbol{\theta}, \mathbf{X}, \mathbf{Z}, \boldsymbol{\Psi}, \sigma^2) &\propto (\sigma^2)^{-\frac{N}{2}} \exp(-\frac{\sum_{i=1}^n \sum_{j=1}^J (y_{ij} - \boldsymbol{\theta}^T \mathbf{x}_{ij} - \mathbf{b}_i^T \mathbf{z}_{ij})^2}{2\sigma^2}) \\ &\propto (\sigma^2)^{-\frac{N}{2}} \exp(-\frac{1}{2\sigma^2} SSR) \end{aligned}$$

where $N = \sum_i^n nj_i$ is the entire sample size (all observations within all clusters). The prior for σ is given above, and therefore we can write the full conditional posterior as:

$$\begin{aligned} p(\sigma^2 | \mathbf{y}, \mathbf{X}, \mathbf{Z}, \boldsymbol{\theta}, \mathbf{b}_i, \boldsymbol{\Psi}) &\propto (\sigma^2)^{-\frac{N}{2}-1} \exp\left(-\frac{1}{2\sigma^2} SSR\right) \\ \sigma^2 | \cdot &\sim IG\left(\frac{N}{2}, \frac{SSR}{2}\right) \end{aligned} \quad (9)$$

Full conditional for $\boldsymbol{\Psi}$ (random effects variance covariance matrix)

Here, we need to write down the posteriors for the different priors we specified. First, let us define the sampling model for the random effects.

$$\begin{aligned} \begin{bmatrix} b_{0i} \\ b_{1i} \end{bmatrix} = \mathbf{b}_i &\sim N(\mathbf{0}, \boldsymbol{\Psi}) \\ p(\mathbf{b}_1, \mathbf{b}_2 | \boldsymbol{\Psi}) &\propto |\boldsymbol{\Psi}|^{-\frac{n}{2}} \exp\left(-\frac{1}{2} \text{tr}(\mathbf{S}_b \boldsymbol{\Psi}^{-1})\right) \end{aligned} \quad (10)$$

where \mathbf{S}_b is $\sum_i \mathbf{b}_i \mathbf{b}_i^T$

- given the inverse-Wishart prior

$$\begin{aligned} p(\boldsymbol{\Psi}) &\propto IW(\nu, \mathbf{S}_0) \\ &\propto |\boldsymbol{\Psi}|^{-\frac{(\nu+k+1)}{2}} \exp\left(-\frac{1}{2} \text{tr}(\mathbf{S}_0 \boldsymbol{\Psi}^{-1})\right) \end{aligned}$$

the full conditional posterior of $\boldsymbol{\Psi}$ is

$$\begin{aligned} p(\boldsymbol{\Psi} | \cdot) &\propto |\boldsymbol{\Psi}|^{-\frac{(\nu+n+k+1)}{2}} \exp\left(-\frac{1}{2} \text{tr}([\mathbf{S}_0 + \mathbf{S}_b] \boldsymbol{\Psi}^{-1})\right) \\ &\propto IW(\nu + n, \mathbf{S}_0 + \mathbf{S}_b) \end{aligned} \quad (11)$$

where $\nu = 2$

- inverse-Wishart *a la* Huang and Wand

$$\begin{aligned} p(\boldsymbol{\Psi} | a_1, a_2) &\propto IW(\nu + k - 1, 2\nu \text{diag}(1/a_1, 1/a_2)), \\ a_k &\propto IG(\eta, 1/A_k^2) \\ p(\boldsymbol{\Psi}) &\propto |\boldsymbol{\Psi}|^{-\frac{(\nu+k-1+1)}{2}} \exp\left(-\frac{1}{2} \text{tr}(2\nu \text{diag}(1/a_1, 1/a_2) \boldsymbol{\Psi}^{-1})\right) \\ &\times \left(\frac{1}{a_1}\right)^{\eta+1} \exp\left(-\frac{1}{A_1^2 a_1}\right) \times \left(\frac{1}{a_2}\right)^{\eta+1} \exp\left(-\frac{1}{A_2^2 a_2}\right) \end{aligned}$$

the full conditional posterior of $\boldsymbol{\Psi}$ is

$$\begin{aligned} p(\boldsymbol{\Psi} | \cdot) &\propto |\boldsymbol{\Psi}|^{-\frac{(\nu+k-1+n+1)}{2}} \exp\left(-\frac{1}{2} \text{tr}([\mathbf{S}_b + 2\nu \text{diag}(1/a_1, 1/a_2)] \boldsymbol{\Psi}^{-1})\right) \\ &\propto IW(\nu + k - 1 + n, \mathbf{S}_b + 2\nu \text{diag}(1/a_1, 1/a_2)) \\ p(a_k | \cdot) &\propto IG\left(\eta(\nu + k), \nu \left(\boldsymbol{\Psi}_{kk}^{-1} + \frac{1}{A_k^2}\right)\right) \end{aligned} \quad (12)$$

where $\eta = \frac{1}{2}$, $\nu = 2$, $k = 2$, and n is the number of clusters (individuals). (For the conditional posterior of a_k refer to Huang and Wand (2013), section 4.2).

- Matrix-F variate

Following section 2.3 in Mulder and Pericchi (2018), instead of working directly with the $\Psi \sim F(\nu, \delta, \mathbf{B})$ we apply the parameter expansion defined above (see section on priors) and model it as $\Psi \sim IW(\delta + k - 1, \mathbf{\Omega})$ with $\mathbf{\Omega} \sim W(\nu, \mathbf{B})$. With this parameter expansion, the conditional priors are:

$$\begin{aligned}\Psi|\mathbf{\Omega} &\sim IW(\delta + k - 1, \mathbf{\Omega}) \\ \mathbf{\Omega}|\Psi &\sim W(\nu + \delta + k - 1, (\Psi^{-1} + \mathbf{B}^{-1})^{-1})\end{aligned}$$

which makes the full conditional posterior of:

$$\begin{aligned}\Psi|\mathbf{\Omega}, . &\sim IW(\delta + k - 1 + n, \mathbf{S}_b + \mathbf{\Omega}) \\ \mathbf{\Omega}|\Psi, . &\sim W(\nu + \delta + k - 1, (\Psi^{-1} + \mathbf{B}^{-1})^{-1})\end{aligned}$$

with parameters as defined above. Given these posteriors, the Gibbs sampler implementation is straightforward.

Notation Conventions

- n number of clusters; i specific cluster
- J number of observations within cluster; j specific observation
- N total number of observations